

A Method and Apparatus for Prefetching Recursive Data Structures

Field of the Invention

This invention addresses the problem of prefetching indirect memory references commonly found in applications employing pointer-based data structures such as trees and hash tables. More specifically, the invention relates to a method for pipelining transactions on these data structures in a way that makes it possible to employ data prefetching into high speed caches closer to the CPU from slow memory. It further specifies a means of scheduling prefetch operations on data so as to improve the throughput of the computer system by overlapping the prefetching of future memory references with the execution of previously cached data.

10 Background of the Invention

Modern microprocessors employ multiple levels of memory of varying speeds to reduce the latency of references to data stored in memory. Memories physically closer to the microprocessor typically operate at speeds much closer to that of the microprocessor, but are constrained in the amount of data they can store at any given point in time. Memories further from the processor tend to consist of large dynamic random access memory (DRAM) that can accommodate a large amount of data and instructions, but introduce an undesirable latency when the instructions or data cannot be found in the primary, secondary, or tertiary caches. Prior art has addressed this memory latency problem by prefetching data and/or instructions into the one or more of the cache memories through explicit or implicit prefetch operations. The prefetch operations do not stall the processor, but allow computation on other data to overlap with the transfer of the prefetch operand from other levels of the memory hierarchy. Prefetch operations require the compiler or the programmer to predict with some degree of accuracy which memory locations will be referenced in the future. For certain mathematical constructs such as arrays and matrices, these memory locations can be computed *a priori*, but the memory reference patterns of the traversals of certain data structures such as linked lists, trees, and hash tables are inherently unpredictable. In a binary tree data structure, for instance, the decision on whether a given traversal should continue down the left or right subtree of a given node may depend on the node itself.

In modern transaction processing systems, database servers, operating systems, and other commercial and engineering applications, information is frequently organized in hash tables and trees. These applications are naturally structured in the form of distinct requests that traverse these data structures, such as the search for records matching a particular social security number. If the index set of a database is maintained in a tree or other pointer-based data structure, lack of temporal and spatial locality results in a high probability that a miss will be incurred at each cache in the memory hierarchy. Each cache miss causes the processor to stall while the referenced value is fetched from lower levels of the memory

hierarchy. Because this is likely to be the case for a significant fraction of the nodes traversed in the data structure, processor utilization will be low.

The inability to reliably predict which node in a linked data structure will be traversed next sufficiently far in advance of such time as the node is used effectively renders prefetching impotent as a means of hiding memory latency in such applications. The invention allows compilers and/or programmers to predict memory references by buffering transactions on the data structures, and then performing multiple traversals simultaneously. By buffering transactions, pointers can be dereferenced in a pipelined manner, thereby making it possible to schedule prefetch operations in a consistent fashion.

Description of Prior Art

Multi-threading and multiple context processors have been described in prior art as a means of hiding memory latency in applications. The context of a thread typically consists of the value of its registers at a given point in time. The scheduling of threads can occur dynamically or via cycle-by-cycle interleaving. Neither approach has proven practical in modern microprocessor designs. Their usefulness is bounded by the context switch time (i.e. the amount of time required to drain the execution pipelines) and the number of contexts that can be supported in hardware. The higher the miss rate of an application, the more contexts must be supported in hardware. Similarly, the longer the memory latency, the more work must be performed by other threads in order to hide memory latency. The more time that expires before a stalled thread is scheduled to execute again, the greater the likelihood that one of the other threads has caused a future operand of the stalled thread to be evacuated from the cache, thereby increasing the miss rater, and so creating a vicious cycle.

Non-blocking loads are similar to software controlled prefetch operations, in that the programmer or compiler attempts to move the register load operation sufficiently far in advance of the first utilization of said register so as to hide a potential cache miss. Non-blocking loads bind a memory operand to a register early in the instruction stream. Early binding has the drawback that it is difficult to maintain program correctness in pointer based codes because loads cannot be moved ahead of a store unless it is certain that they are to different memory locations. Memory disambiguation is a difficult problem for compilers to solve, especially in pointer-based codes.

Prior art has addressed prefetching data structures with regular access patterns such as arrays and matrices. Prior attempts to prefetch linked data structures have been restricted to transactions on those data structures in which the traversal path is largely predictable, such as the traversal of a linked list or the post-order traversal of a tree. The invention described herein addresses the problem of prefetching in systems in which the traversal path is not known *a priori*, such as hash table lookup and tree search requests. Both of these traversals are frequently found in database applications, operating systems, engineering codes, and transaction processing systems.

Summary of the Invention

The present invention significantly increases the cache hit rates of many important data structure traversals, and thereby the potential throughput of the computer system and application in which it is employed. The invention is applicable to those data structure accesses in which the traversal path is dynamically determined. The invention does this by aggregating traversal requests and then pipelining the traversal of aggregated requests on the data structure. Once enough traversal requests have been accumulated so that most of the memory latency can be hidden by prefetching the accumulated requests, the data structure is traversed by performing software pipelining on some or all of the accumulated requests. As requests are completed and retired from the set of requests that are being traversed, additional accumulated requests are added to that set. This process is repeated until either an upper threshold of processed requests or a lower threshold of residual accumulated requests has been reached. At that point, the traversal results may be processed.

Brief Description of the Drawings

Figure 1 is a high level representation of a typical transaction processing system. As each request is received, actions are performed based on the event type.

Figure 2 is a high level representation of a transaction processing system that employs the invention in order to accumulate events that require traversal of the same type of data structure. Events are buffered until a threshold has been reached, at which point software pipelined prefetching is performed across some or all of the accumulated events.

Figure 3 illustrates an extension to the method of *Figure 2* by buffering the results of the event processing, in this case database searches, for subsequent simultaneous traversal.

Figure 4 illustrates the event buffering process on a tree in which event processing is deferred until at least 4 search requests have been received. The state of the event buffer *S* is illustrated after each of the requests *R*₀, *R*₁, *R*₂, and *R*₃, have been received.

Figure 5 illustrates the state of the system as the four search requests traverse the data structure. The state of the event buffer *S* is illustrated after matching keys for a search request have been found, first for *R*₂, and then *R*₀. Results are stored by swapping the satisfied request with the last unsatisfied event in the buffer. Once a completion threshold has been reached, in this case 2, the search is discontinued until the accumulated results have been processed.

Figure 6 illustrates the state of the system after two additional requests, *R*₄, and *R*₅, have been buffered. The search continues as in *Figure 5*, except that, this time, the number of nodes with matching search requests is allowed to go above the completion threshold. The matches may exceed the completion threshold because all requests for which a prefetch has been issued are allowed complete processing of the prefetched node, thereby preventing memory references from subsequent result processing, operating system activity, and other processes from displacing the prefetched data.

Figure 8 shows the structure of a software implementation of the invention.

Figure 9 illustrates the process of accumulating K requests on accumulation queue AQ for software pipelined traversals of data structure S , where K is the startup threshold. Accumulated results are returned from result queue RQ .

5 *Figure 10* shows a pseudo-code example of an initial call to a recursive binary tree search.

Figure 11 provides a pseudo-code description of the recursive component of a pipelined search of a binary tree. The request is added to the result queue when the current node is NIL, indicating that a node with a matching key does not exist in the tree, or the key of the current node matches the requested key.

10 *Figure 12* shows memory alignment of hash table slots. The dark shaded areas of the unaligned hash slots indicate slots that have the potential to generate two misses; the light shaded slots miss at most once per reference. Packed slots are smaller because the pointers to the next hash table entry are eliminated from all but the last element in a packed slot.

Figure 13 shows how heterogeneous hash table data structures (a) are represented as homogeneous structures (b), thereby eliminating at least one memory reference.

15 *Figure 14* shows how some potential cache misses in (a) can be eliminated by padding the data structure so that elements always fall on hash line boundaries in (b).

Figure 15 Hash table packing. A homogeneous hash table structure (a) may be represented as a packed structure (b), which can be rebalanced to make the table less sparse as in (c), (d), (e), or others.

20 *Figure 16* illustrates the architecture of a Transaction Buffer, used to assist in accumulation and processing of traversal requests on a data structure.

Detailed Description

Prefetching pointer-based data structures is much more difficult than prefetching data structures with regular access patterns. In order to prefetch array based data structures, Klaiber and Levy¹ proposed using software pipelining – a method of issuing a prefetch request during one loop iteration for a memory operand that would be used in a future iteration. For example, during loop iteration j in which an array $X[j]$ is processed, a prefetch request is issued for the operand $X[j+d]$, where d is the number of loop iterations required to hide the memory latency of a cache miss. The problem with this loop scheduling technique, prior to the introduction of this invention, is that it could not be applied to pointer-based data structures. A concurrently submitted application addresses this problem for data structures in which the traversal path is predefined, such as linked list traversal and post-order traversal of a tree. This invention addresses the problem for data structure traversals in which the traversal path is dynamically determined, such as in hash table lookup and binary tree search traversals. The application of the invention is then illustrated by means of binary search trees and hash table lookup.

The invention consists of the following method. Step 1 is to homogenize the data structure(s) to be traversed, where applicable. This process is described for open hash tables below and illustrated in figure 13, and can be applied to adjacency lists commonly used to represent graphs, and other data structures in order to remove an unnecessary level of indirection. Step 2 is to align the data structure on cache line boundaries, where applicable, as described below for hash tables and illustrated in figure 14. Alignment is performed for each element in the data structure, i.e. for each node of heap allocated storage in a linked list or tree. In an array, this may mean introducing pad space, such as is described for hash tables below. In heap allocated nodes, it may mean employing a special heap management interface that aligns blocks of memory on cache line boundaries. Step 3 packs data structures into cache line size components, where applicable, as described for hash tables below and in figure 15. The data structure traversal is constructed so that traversal is first performed over a group of data structures that have been packed together, and then over potentially multiple groups of packed structures. Step 3 is applicable primarily to nodes which exhibit temporal locality, such as nodes in a linked list (for instance, in an open hash table or adjacency list representation of a graph). Step 4 is to buffer events on a given data structure until a critical number of events has been reached, as described below for trees and hash tables, and as illustrated in figures 2 through 4 and figure 9. The exact number of events to buffer can be determined via the compiler, experimentally, or a combination of both. Step 4 can be implemented by returning to the caller when the number of buffered requests is below a given threshold, unless immediate traversal is requested by another component of the system. Step 5 is to traverse the data structure for which the events have been buffered in a pipelined manner, issuing prefetch requests for some or all of the buffered events. Step 5 is illustrated for select examples in figures 5, 6, 10, and 11 and is described for search tree traversals

¹ Klaiber and H.M. Levy, An Architecture for Software-Controlled Data Prefetching, Proceedings of the 18th International Symposium on Computer Architecture 1991, pp. 43-53.

and hash tables below. The traversal results, such as when a node for which a matching key has been found, can either be processed immediately, or be stored in a buffer as well. The prefetch distance may be determined experimentally by the programmer, computed using prior art, or determined by the compiler. There are three parameters that control the traversal: a startup threshold, a completion threshold, and the pipeline depth. The startup threshold is the number of events that are buffered before traversal of the data structure or structures is allowed to commence. The completion threshold is the number of traversals that are completed (for example, when a matching key has been found) before no additional events are processed from the accumulated event buffer. Step 6 passes the results to the next stage in the surrounding system. The overall process is illustrated in figure 8. The buffers can be implemented in hardware or in software.

The data structures and traversal algorithms addressed in the concurrently submitted application have a common feature: only a single traversal path is taken through the data structure. Data dependencies may affect *whether* the path is taken to completion, which does not materially affect the choice of prefetch targets. The property that the path through a data structure is independent of the values of the nodes within the data structure makes it possible to modify the data structure so that the necessary parallelism to support software pipelining can be exposed. This condition does not hold for tree and hash table searches. In this application I discuss a method of aggregating temporally distributed data structure traversals in order to support software pipelined prefetching, which I refer to as *temporal restructuring*.

d	Prefetch distance
l	Prefetch latency
n	Node on a traversal path.
n_r	Root or startup node
n_t	Termination node of traversal
K	Startup threshold
P	Pipeline depth
Z	Completion threshold
AQ	Accumulation buffer
IQ	Prefetch issued buffer
RQ	Result buffer
R	Request $\langle k, r, a \rangle$
k	Search key
a	Node address
r	Request id

Table 1: definition of variables

Search paths through a tree are not generally predictable. Consider a path P from the root node n_r to a termination node n_t , $P = n_r, n_1, \dots, n_t$. In the case of a binary search tree, the value of n_t depends on the value of the key field of n_{t-1} .

If both the left and right node of a tree are always prefetched, then one prefetch target will usually have been prefetched in vain. Prefetching can only be effective if the prefetch address is identifiable far enough in advance so that it can be prefetched into near memory by the time it is first referenced. Thus even if both children are prefetched, a single pass of the inner loop of the search below does not require
 5 enough cycles to hide any significant latency:

```

Node SEARCH(Tree tree, Key k)
begin
  Node n;
  10  n ← tree.root;
  while ( node != NULL ) do
    if ( k = n.key ) then return n; endif
    if ( k < n.key ) then
      15  n ← n.left;
    else
      n ← n.right;
    endif
  end while
  return NIL;
  20  end

```

Consequently, the ability to prefetch only the children of the current node is unlikely to provide sufficient computation between the time the prefetch is issued and the time it arrives. In general, software pipelining schedules prefetch operations $d \geq \lceil l/s \rceil$ loop iterations ahead in order to completely hide latency, where s is
 25 the execution time of the shortest path through the loop and l is the prefetch latency. If the prefetch distance d is small, and the tree has been mapped to an array, it may be possible to employ greedy prefetching of the entire sub-tree of depth d . I refer to this subtree at node n_i as the *prefetch subtree* of n_i . For the root node, the entire subtree of 2^{d-1} nodes would have to be prefetched, of which all but d are prefetched in vain. For each of the subsequent $p-d-1$ nodes on the path, 2^{d-1} nodes are prefetched, resulting
 30 in $(p-d-1) \times 2^{d-1} - 1$ extraneous prefetches. The last $d-1$ nodes on P correspond to the epilogue in traditional software pipelined prefetching, requiring no additional prefetch commands. These numbers may actually be optimistic, since they assume that the application can avoid prefetching the entire subtree of 2^{d-1} nodes at each node in the path, issuing prefetches only for the newly discovered 2^{d-1} leaf nodes of the prefetch subtree. It is obviously not desirable to prefetch up to 2^{d-1} nodes when only 1 is required at each node
 35 along the path.

While a single traversal of the tree does not provide sufficient opportunity to exploit software pipelining, I show how temporally scattered, independent search requests can be aggregated so that software pipelining can be applied across multiple requests. The premise behind the approach is that a
 40 single unit of work performed on a given data structure may not provide sufficient opportunity to hide the latency via software pipelining, so work is allowed to accumulate until a threshold is reached, or a request for immediate resumption forces work to proceed. I refer to this process of aggregating and collectively

processing independent data structure traversals as *temporal restructuring*. Figure 9 shows high-level pseudo-code for handling a traversal request. Postponement is implemented by returning the special token POSTPONE from the search routine, which can be considered a request to submit additional work. The software pipeline processes accumulated work, storing results in the queue RQ.

5

In an online transaction processing environment, for instance, multiple temporally proximate transactions can be grouped for simultaneous traversal of the data structure. The amount of time that any particular search can be postponed in a transaction processing system may be limited by system response time requirements. Since the number of search requests that must be accumulated in order to ensure a software pipeline depth adequate to effectively hide memory latency is relatively small (in the tens of requests), this should not be an issue in a high throughput system. A system with real-time constraints must be able to ensure completion even when the system is not very busy. Since the number of search requests can be adjusted dynamically, the startup threshold, K in figure 9, can be reduced to a single request. In general, systems can raise the threshold in order to improve throughput and lower it to improve system response time.

15

The general structure of the accumulation process is illustrated in figure 8. A search request consists of a *transaction descriptor* $\langle k, r, a_i \rangle$, where k is the search key associated with the request data structure identifier r , and a_i is the initial prefetch target address, such as the address of the root node of a search tree or the initial hash bucket in the bucket chain of a hash table. One or more *prefetch descriptors* are associated with each data structure. A *prefetch descriptor* stores the invariants among accumulated requests, characterizing the prefetch target by the pipeline depth D , the *startup threshold* K , the *completion threshold* Z , the number of bytes to prefetch at each request, and a small buffer for several words of application-defined data.

20

25

Search requests are accumulated in AQ , the *accumulation queue*. When the number of elements in the queue reaches the startup threshold, K , then D search requests are dequeued from the accumulation queue. The address portion of each request is submitted to the prefetch hardware along with the prefetch parameters from the prefetch descriptor, and the request is enqueued on the *prefetch issued queue*. This sequence of actions corresponds to the prologue of software-controlled prefetching.

30

The accumulation process for a binary search tree is illustrated in figure 4. In this example, $D = 2$ and $K = 4$. The accumulation queue is shown after the requests R_0, \dots, R_3 have been submitted. Once R_3 has been submitted, the completion threshold is reached, and all four requests point at the root node. The application dequeues the *active request*, $\langle k, r, a_c \rangle$, from the prefetch issued queue and processes the data corresponding to the *current address* a_c . If the requested key matches the key at the current address, then the active request is enqueued on the *result queue*, and a prefetch is issued for the next request on the

35

accumulation queue. This request is then added to the prefetch issued queue. If the keys do not match, then a prefetch is issued for the next address a_n in the data structure (such as a child pointer). The new address a_n replaces a_c in the active request transaction descriptor, and the new transaction descriptor $\langle k, r, a_n \rangle$ is added to the prefetch issued queue.

5

When the application is ready to process a search result, it extracts a search result descriptor $\langle k, r, a_x \rangle$ from the *result queue*, where a_x is the address of the node containing k . Applications that perform searches typically return a value of NIL to indicate that no matching key was found. Thus, when no matching key is found, a_x can be set to NIL. The application can force the search to proceed, even when the number of accumulated search requests falls below K , by submitting a reserved *completion descriptor* for accumulation.

10

If the number of cycles that is required to process the result is small, it may make sense to process each result immediately, rather than adding it to the result queue for later processing. It is not generally desirable to process results right away, since result processing may increase the amount of time spent at a single beat of the software pipeline. Increasing the amount of processing spent at one beat increases the danger that previously prefetched memory locations will again be displaced from the cache. If processing the result requires any I/O, for instance, the processor is likely to suspend the current process and perform other work. It is quite possible that all outstanding prefetches will be overwritten in the cache before the process that issued them is scheduled to run again. In the worst case, it is scheduled to run in a different CPU².

15

20

² For example, Xia found that some operating system activity involves clearing large buffers, which invalidates a large number of cache entries.

Binary Search Trees

The method can be demonstrated by applying it to a binary tree search. In order for the technique to be applicable, multiple independent search requests must be available for processing. To provide a context for a set of search tree traversals, consider a generic processor-bound client/server application that processes account debit requests. A non-pipelined version is illustrated in high-level pseudo-code below:

```

GENERIC-SERVER()
begin
    loop forever
        request ← GetNextRequest();
        case request.type of
        DEBIT: begin
            FindAcct( request );
            display ← OtherWork( request );
            Reply( display );
        end
        OTHER: ...
        end case
    end loop
end Server

```

search key, and a pointer to an account record that is used to hold the search result. This query data structure strongly associates the query result with the answer, making it easier to support multiple outstanding queries. A viable alternative implementation might have the search routine return a pointer to the caller. A prerequisite of temporal restructuring is the ability to strongly associate a request with a request result, so that downstream code can work on a different request than that submitted to the search routine. Rather than cluttering the examples with implementation details of the straightforward process of associating requests with results, the example starts with an implementation in which the search result is bound to the search request as a pre-existing condition. Thus, the server searches a database for an account record corresponding to the search key, and the account pointer is initially set to NIL. The application invokes a search routine *FindAcct* that traverses the index set organized as a tree. *FindAcct* updates the account pointer to the matching account record, if there is one, and NIL otherwise. Based on the account information, a response to the request is generated in other subsystems, indicated by the call to *OtherWork*. The response is returned to the client application via the call to *Reply*.

A version of the server that processes DEBIT requests in a pipelined manner is illustrated below:

```

PIPELINED-SERVER()
begin
  loop forever
    request ← CheckNextRequest();
    case request.type of
      DEBIT: begin
        qPipeSubmit (WorkQ,request) ;
        result ← qPipeExtract (ResultQ) ;
        if result ≠ NIL then
          display ← OtherWork (result) ;
          Reply (display) ;
        endif
      OTHER: ...
    end case
  end loop
end

```

The search tree traversals in this version are performed as part of *qPipeSubmit*, but only once the number of requests in the pipeline has reached K . When fewer than K requests occupy the pipeline, no search requests are processed, and *qPipeExtract* returns the reserved address `NONE_AVAILABLE`.

- 5 Otherwise, *qPipeExtract* returns the first request for which a search result is available.

- In an online transaction processing (OLTP) environment, *GetNextRequest* is a blocking call, stalling the server thread or process until another request becomes available. *CheckNextRequest* is a modified version of *GetNextRequest* which returns a synthetic DEBIT request containing a completion descriptor that forces any pending accumulated requests to complete if the result queue is empty. If the result queue is not empty, the application extracts a completed request as before, albeit without enqueueing a new request. Thus the server stalls only when all accumulation queues and result queues are empty, which avoids delaying replies when the request arrival rate is low. Although the system would not achieve maximal efficiency unless the pipeline is filled, the decrease in the arrival rate indicates that the system is otherwise idle, and wasted cycles less precious. In an offline data processing environment, completion is forced after the last request has been submitted.

- If all requested keys are represented in the tree, calls to *qPipeSubmit* simply return until the number of requests submitted to *qPipeSubmit* reaches K . Once K requests have accumulated, prefetches are submitted for the first D requests in the accumulation queue. Each time a prefetch is submitted, the corresponding request is removed from the accumulation queue and added to the prefetch issued queue. This sequence of events constitutes the prologue. Once the prologue has completed, the head of the prefetch issued queue is removed and the corresponding node is processed. If the keys of the request and the node match, then the node address is saved and the descriptor is added to the result queue.

Otherwise, the current descriptor is updated with the appropriate child pointer. At this point, the implementor or compiler has several choices of prefetch strategies:

1. Prefetch the child pointer and add the current request to the end of the accumulation queue. This approach maximizes the prefetch distance, the available distance between the time the prefetch is issued and the corresponding request is again processed. Increasing the prefetch distance beyond the minimum needed to hide memory latency also increases the risk of additional cache conflicts.
2. Issue a prefetch for the next request on the accumulation queue and move that request to the end of the prefetch issued queue. Then move the current request to the end of the accumulation queue without prefetching its child pointer. This process ensures that each of the accumulated requests is processed in round-robin order. Note that round robin scheduling does not guarantee any particular completion order among accumulated requests, since search requests may complete at any point in the traversal of a tree.
3. Prefetch the child pointer and add the current request to the end of the prefetch issued queue, ensuring that requests are processed in approximately first come first served order. Each request remains in the prefetch issued queue until it completes, guaranteeing a larger percentage of processing time to queries that have reached the prefetched issued queue ($1/p$ instead of $1/(p+a)$, where p is the number of elements in the prefetched issued queue and a is the number of elements in the accumulation queue). Once again, this approach does not guarantee a completion order. If the system has been appropriately tuned, the processing delay provided by the queue of length D should be sufficient to hide any memory latency.

The latter two options have similar interference and throughput characteristics. In any case, the head of the prefetch issued queue is removed to replace the current request, and the process repeats itself until no more requests occupy the prefetch issued queue.

If an address for which a prefetch request has been issued is referenced before it has arrived in the cache, the CPU stalls until the prefetch for the corresponding line completes. The cache hardware checks if there is an outstanding prefetch request for the line, in which case it does not issue a second request. Consequently, a reference to an address for which a prefetch is in progress incurs only a partial miss penalty. Cache misses always bypass pending prefetches in the reference system, so that a cache miss never has to wait for the prefetch issue queue to empty before it is submitted to the memory hierarchy. A prefetched address may be evicted before it is referenced for the first time, either by another prefetched

cache line or by another data reference. In this case the CPU stalls until the corresponding line is brought into the cache through the normal cache miss mechanism, incurring the full miss penalty.

As the number of elements in the result queue increases to the point where fewer than D requests remain in the accumulation and prefetch issued queues, not enough work is left to hide the latency. Consider the point in the search process where all but a single search request has been resolved. The request is dequeued from the issued queue, its node pointer updated, and a prefetch issued. It is added to the prefetch issued queue, and almost immediately dequeued again. A prefetch that goes all the way to memory is unlikely to have completed by the time this last search request has been dequeued again, causing the processor to stall.

To avert this problem, I employ a completion threshold Z . As long as the combined number of requests remaining in the issued queue and the accumulation queue remains above Z , a prefetch request is issued for the child pointer, and its descriptor is added to the end of the prefetch issued queue. Once the completion threshold Z has been reached, the current descriptor is added to the accumulation queue instead, without issuing a prefetch request. Inserting the descriptor at the head of the queue, instead of the tail allows search requests that have been waiting longest complete sooner. The remaining elements in the prefetch issued queue are then processed, so that the prefetch issued queue is empty by the time the application exits the epilogue. When there is little danger that the temporarily abandoned prefetch requests will be evicted before the corresponding search requests are resumed, then there is no need for an epilogue. This information is not generally predictable at compile time, and since the epilogue has only a moderate impact on the instruction cache footprint of the application, it is generally included. This process may bring the actual number of remaining requests below Z , since some of requests in the issued queue may move to the result queue. All other requests are available on the result queue, and will be dequeued and returned by repeated calls to the result extraction routine, *qPipeExtract*, until the result queue is empty. The PIPELINED-SERVER above shows *qPipeExtract* being used to process results from the software pipeline.

Intuitively, it would appear that D is a natural choice for the value of the completion threshold Z , with K some small multiple of D . Yet, in experiments where Z was varied from 0 to D and K was kept constant, a relative performance declined notably well before the completion threshold Z reached the pipeline depth D . This is a consequence of the fact that the amount of work performed per set of accumulated traversals decreases as Z approaches D . The amount of work performed each time traversal is triggered is a function of $K - Z$. For a fixed value of K , the amount of work performed at each traversal decreases as Z increases. If the amount of work performed with each traversal is decreased, then more traversals are required to accomplish the same total amount of work. For instance, if there are 1000 requests and $D = 22$, $Z = 22$, and $K = 32$, then only 10 requests are completed per traversal, requiring 100

traversals. Part of the startup cost of each traversal is a function of D . Another portion of the startup cost is bringing the working set that is not prefetched, such as instructions that have been evicted between traversals, into cache. The startup overhead is incurred 100 times for 1000 requests. If Z is reduced to 12, then the startup cost is incurred only half as often, although each traversal will have to endure at least

5 partial latency due to a partially empty software pipeline. The application effectively trades off some latency for a reduction in startup overhead.

Because search traversals of data structures typically perform very little work at each node, the optimal pipeline depth can be quite long. A tree search achieved optimal performance at a pipeline depth

10 of 32, while the optimal pipeline depth for Quicksort was only four. Quicksort performs significantly more work per iteration. Response time constraints and interference effects may limit the practical size of the accumulation queue. If queue management is supported by hardware, hardware constraints will further curtail the total number of outstanding requests that can be accommodated.

15 The startup threshold used to begin a round of searching is adjusted so that most of the latency can be hidden most of the time without violating other system requirements such as service response time. For software-managed queues, an accumulation scheme that attempts to accumulate too many requests may also introduce self-interference, since the queues also increase the cache footprint of the application.

20 When the first round of searches is triggered, the node set will contain only pointers to the root of the tree. This will result in multiple prefetch instructions for the same node. The memory hierarchy keeps track of pending prefetch requests, ensuring that only one memory request is outstanding to the same cache line at any given time. Consequently, multiple prefetch requests to the root node do not generate additional memory traffic. After the initial cold start, a diverse set of partially completed requests will populate the

25 set of nodes in the accumulation queue. Some of the search requests in the accumulation queue may be the result of a traversal that reached the completion threshold, and thus refer to arbitrary nodes within the tree. Maintaining this state may impose some restrictions on insertion of nodes into the tree and deletion of nodes from the tree, as is discussed below.

30 If program semantics allow search requests to intermingle with insertion and deletion requests, then software pipelining introduces some new timing issues, especially in the presence of completion thresholds. When node insertion and deletion are supported, then the fact that there may be search requests already in progress may impact the outcome of the searches. Consider a search tree undergoing insertion of a node with key k_n , and that key does not exist in the tree at the time of the insertion. There may be an

35 outstanding search for k_n , invoked prior to the insertion. Had the request been processed immediately, it may have returned a NIL pointer. If insertion requests are processed concurrently with the search requests, then the insertion may complete before the outstanding search request for k_n is processed. If the node is

inserted at a leaf and the tree undergoes no other changes, then the next time the number of accumulated searches reaches the startup threshold, k_n will find the freshly inserted node. In doing so, it will return a node that did not exist at the time the search request was first made. This clearly changes the semantics of the program. In some cases, this may be of no particular consequence. For instance, consider a banking system that allows creation of accounts on a live online transaction processing system. If the accounts are indexed by social security number via a tree, then it is very unlikely that a search request for an account will precede the creation of that account. Even if this unlikely event comes to pass, a positive outcome of a search is not necessarily regarded as an error. Since temporal structuring is aimed at hiding memory latency, the time spans can be measured in fractions of seconds - for applications involving human interfaces this is practically simultaneous.

For balanced binary tree schemes, where the table is modified with each insertion, it may be prudent to force completion of extant search requests. Consider the insertion of node *A* in figure 7. The resulting rotation moves node *B* to the root node. This means that all pending search requests pointing to *C* will never find *B*, since search requests always traverse downwards. The problem with deletion is more obvious, since a node that is pointed to by an accumulated search request for which processing is not complete, cannot be safely deleted without affecting the search operation. If any search request in the accumulation queue points to the deleted node, then the outcome is almost guaranteed to be a disaster if the memory occupied by the deleted node is reused for some entirely different purpose. Since deletions from the tree between search requests may invalidate accumulated search requests, either deletion can be postponed until all accumulated searches have completed, or pending search requests can be forced to complete prior to deletion.

Recursion

Recursion is often a natural way to express an algorithm. Recursive tree searches can be performed in a manner similar to the loop-based tree search described above. Figures 10 and 11 illustrate a recursive version of a search. The code shows the explicit management of the request and result queues to illustrate how a software-only approach can be employed. Once enough requests have been accumulated, the recursive search, TREE-DELAYED-SEARCH of figure 10 is called to execute the prologue. All requests recurse together. That is, each recursive call to TREE-RECURSE advances every pending request one node in the tree. The recursion is allowed to unravel once a sufficient number of search requests have been satisfied so that the number of outstanding search requests no longer meets the lower bound criterion on the software pipeline depth. The binary tree search in the example is tail-recursive, so that the state of the search requests does not need to be maintained on the stack. Consequently, allowing the recursion to unravel has no negative repercussions.

For applications that rely on maintaining the state of the stack variables from prior procedure invocations, allowing the recursion to unravel could prove more of a problem. In these cases, all searches in the pipeline can be allowed to complete, without regard for the completion threshold, and at the expense of more memory stalls.

A more general version of the tree traversal algorithm would have to place all nodes and keys onto the stack, considerably increasing the amount of stack space required to complete the search, and thus the data cache footprint of the application.

Hash Tables

Hash tables and other data structures with short pointer chains pose a particular challenge to prefetching. The problem is two-fold: short pointer chains do not provide much to prefetch, and the amount of work performed at each iteration in the process of prefetching them is negligible, therefore actually requiring a significant prefetch distance in order to hide memory latency. This patent includes several methods to cope with this problem.

The hash table data structure is modified so that, instead of storing a pointer to a list of hash buckets, each hash table slot contains the first bucket of each chain in the hash table directly. Empty entries are indicated by an invalid key. If there are no invalid keys, an empty entry can be indicated via a reserved address in the pointer to the next hash bucket. This optimization serves two purposes. First, it eliminates one unnecessary level of indirection by allowing the hash function to directly supply the address of the first bucket in the chain. Second, it has the effect of homogenizing prefetch targets. Homogeneous prefetch targets eliminate the need for separate code to prefetch the initial hash table entry. This has the effect of increasing the size of each hash table slot, which should only prove disadvantageous if there are a preponderance of empty hash table slots. If the hash table is fully populated, then I've actually reduced the memory requirements by the size of the hash table.

The homogenized hash table can be subjected to several locality optimizations. An obvious means of eliminating cache misses is to ensure that each hash bucket is aligned on a cache boundary. Hash buckets in the benchmarks that I used to evaluate the efficacy of the approach consist of a key, a pointer to the next element on the hash table, and a pointer to a data record, for a total of 12 bytes. If each L1 data cache line supports 16 bytes, half of the hash table slots will span two L1 cache lines, and every third hash table slot will span two 32 byte L2 cache lines, as illustrated in figure 13. The problem can be eliminated by padding each entry so that the subsequent array element is aligned on a cache line boundary. Hash buckets in the hash chain are similarly aligned. This approach does not require any fundamental change in the hash lookup code. Both prefetching and non-prefetching versions can expect to benefit from alignment.

Cache line sizes in modern microprocessors are 32 bytes or more in size for primary caches, and 64 bytes or more for secondary caches. Large line sizes are an invitation to pack more data into each cache line. Additional performance benefits can be derived by packing as much of a bucket chain into each cache line as possible. This approach appears attractive when the hash chains contain more than a single element. Note that long hash chains run contrary to the philosophy of hashing. Adjacency lists employed by many graph algorithms, on the other hand, may maintain an arbitrarily long list of references to adjacent nodes in the graph.

Alignment and homogenization help reduce the number of cache misses incurred in a hash table lookup. When a hash chain achieves any significant length, the hash chain can be packed into a buffer that fits into a cache line. The buffer is structured as an array of hash chain elements followed by a pointer to the next buffer. For a buffer containing n elements, the pointer to the next hash element can be eliminated for the first $n-1$ elements, allowing more hash chain entries to be accommodated in each buffer. A reserved key can be used to indicate the end of the array, or the pad word can be used to hold the number of valid hash chain entries in the array. The last word in the buffer is used to hold the address of the next buffer, allowing for the possibility that the length of the hash chain may exceed the number of elements that can be accommodated in a single buffer. In a sense, the implicit prefetch inherent to large cache lines is being employed for buckets that share a cache line. Explicit prefetching can be applied to prefetch each packed buffer, thereby increasing the likelihood that a cache line will be available if the number of collisions should exceed the capacity of a single packed hash line, with the added benefit that each prefetch operation can actually prefetch up to n hash chain elements.

The number of hash collisions per bucket can be expected to remain small under ordinary hashing conditions and given a good hash function. Aligned hash chains have the disadvantage that the minimum size of a hash chain is relatively large. A packed buffer for an 8 word cache line contains 3 entries of 2 words each for the example, in addition to a word of padding and a pointer to the next hash bucket. Figures 15 (b) through 15 (e) illustrate the configuration of a packed hash table. As few as one out of three hash chain entries may actually store hash elements when the number of collisions is small. Consequently, packing may lead to a fair amount of wasted memory, especially in a homogenized hash table. Based on the assumption that hash collisions are cheap if they fall onto the same cache line, I reduced the number of hash table slots for the packing transformation, modifying the hash function accordingly. While the number of collisions at each slot increases, the cost of resolving a collision decreases. If the average number of entries per hash slot is originally 1, then the only benefit afforded by the approach is a slight increase in the probability that a hash bucket may be reused across multiple instantiations. For the benchmarks, hash buckets were packed into an array, with each array aligned on an L2 cache line boundary.

Experimental results show that packed hashing afforded similar benefit to alignment. Packing improved hash lookup performance by 4% to 17%, while alignment, alone, improved hash lookup performance by 2% to 16% when the average hash chain contained a single element. This significant improvement indicates that hash elements that span multiple cache lines have a significant negative impact on hash lookup performance. When the average hash chain length increases to 1.5, alignment affords an 8% to 10% performance improvement. Temporal restructuring, when applied to hash tables without specialized hardware support, did not perform well, since the overhead is amortized over few memory references. Performance improved from 4% to 14%, depending on hash interference assumptions between requests. Combining alignment and prefetching did not significantly improve the performance, showing a 12% to 20% performance improvement over the non-prefetching implementation. Experiments showed that hardware buffers, as illustrated in figure 16, significantly improve the performance of the system beyond that of the original prefetch scheme.

Hardware Buffer

The mechanism for buffering transactions described thus far employ buffers allocated from general memory. System throughput can be significantly improved by providing this buffer in hardware, along with a few operations on the buffer.

Figure 17 illustrates a hardware implementation of the temporal restructuring mechanism described in figure 9. The queues are implemented in hardware. Each queue is accessed via an independent device port or register. Ports are modeled as memory mapped devices, with a separate address for each field of a tuple. An internal 64-bit control register holds the prefetch descriptor; a separate 32-bit word describes the state of the prefetch queues. The control register holds the values of the programmer- or compiler-supplied software pipeline depth and the completion and startup thresholds. The expiration field is reserved for a timeout value, so that an application can specify the maximum time that is allowed to elapse between the time an element is added to an empty queue and a result should be processed. Control bits allow the system to specify the handling of timer expiration events, which may include forcing completion at the next request or generating a hardware exception. The prefetch target descriptor is used by the application to indicate the location of the prefetch targets relative to the prefetch address. In order to support prefetching of selected fields within a given data structure, the architecture supports selected prefetching of memory locations on a word basis within an address range. Two interpretations of the prefetch target descriptor are available, depending on the value of the mode bit, M : A bit value of 1 at offset w in the word mask indicates that the word at an offset of w words from the specified prefetch address is a prefetch target. For a 4 word (32 byte) line size, the 32 bit word mask makes it possible to efficiently select prefetch targets from among 4 adjacent cache lines with a single prefetch request. The resolution of a single word makes it possible for the prefetch hardware to determine when any multi-word prefetch target falls on one or two cache lines, since prefetch addresses need not be aligned on cache line

boundaries. For instance, a prefetch target with an address A_0 and a word mask of $0xA0000000$. The word mask indicates that the first two words following A_0 are to be prefetched. A single cache line prefetch is required when $A_0 \bmod 32 = 0$. However, an address A_1 , where $A_1 \bmod 32 = 28$, requires two cache lines to be prefetched. Alternatively, each of four 8-bit fields specifies a word offset and a word count of 4 bits each. This makes it possible to specify prefetch targets at four separate offsets from the prefetch base address of up to 16 words each, assuming offsets are cumulative. Other combinations are possible. Additional options require adding bits to the mode field. This flexibility in specifying a prefetch target makes it possible to address the situation where multiple data fields used by the application to process each node fall onto different cache lines, and the data fields do not share a cache line with the key field used to traverse the data structure.

The application writes request tuples to the accumulation queue port, represented by register A of figure 17. Once all components of the request tuple have been received, the request is enqueued on the accumulation queue. Completed requests may be stored on a hardware result queue. If any completed requests occupy the result queue, then the runtime system dequeues one each time a new request is submitted, thus ensuring that the result queue never overflows. Similar to its software counterpart, the hardware result queue indicates that no request is available on the result queue by placing a special value in the result extraction register, X . Each of the hardware queues behaves similar to its software analog. An update of the *next address* register, N , updates the prefetch address field of the active request and places it on the prefetch issued queue. Enqueuing a result causes a request to be moved from the accumulation queue to the issued queue. The prefetch address generator dequeues the head of the issued queue. Based on the prefetch target descriptor, it then generates a sequence of prefetch operations. The request is immediately placed on the prefetch issued queue, rather than waiting for prefetch address generation to complete for the request. In this manner, the application partially hides the latency when the number of available requests to be processed falls beneath the minimum software pipeline depth. I found that this eliminated up to 16% of the read stall cycles in some instances.

The application has the option of placing the result in the result queue via the result register, R . The result queue is present to allow the application to maintain software pipeline semantics. The presence of the result queue does not prevent the application from processing a result immediately, in which case it may be neither necessary, nor desirable, to add the result to the result queue. A system library provides the necessary interfaces to the prefetch unit. Table2 provides an overview of the interface macros provided to support temporal restructuring in hardware.

accumulate	writes the request to the accumulation queue port. A NULL request indicates to the prefetch unit that the pipeline is to be forced.
-------------------	---